# RNA-seq+ - Analysis

**Vojtěch Bystrý**

**16. December 2019**

# NGS experiments

# NGS experiments

```
                        ┌─────────────────┐      ┌─────────────────┐
                        │                 │      │   Recognize     │
                   ┌───▶│ DNA sequencing  │─────▶│ differences from│
                   │    │                 │      │    "normal"     │
┌──────────────┐   │    └─────────────────┘      └─────────────────┘
│ Next Generation│  │
│  Sequencing   │──┤
└──────────────┘   │    ┌─────────────────┐      ┌─────────────────┐
                   │    │  RNAseq, Chip-  │      │    Counting     │
                   └───▶│ seq, ATAC-seq,..│─────▶│    elements     │
                        └─────────────────┘      └─────────────────┘
```

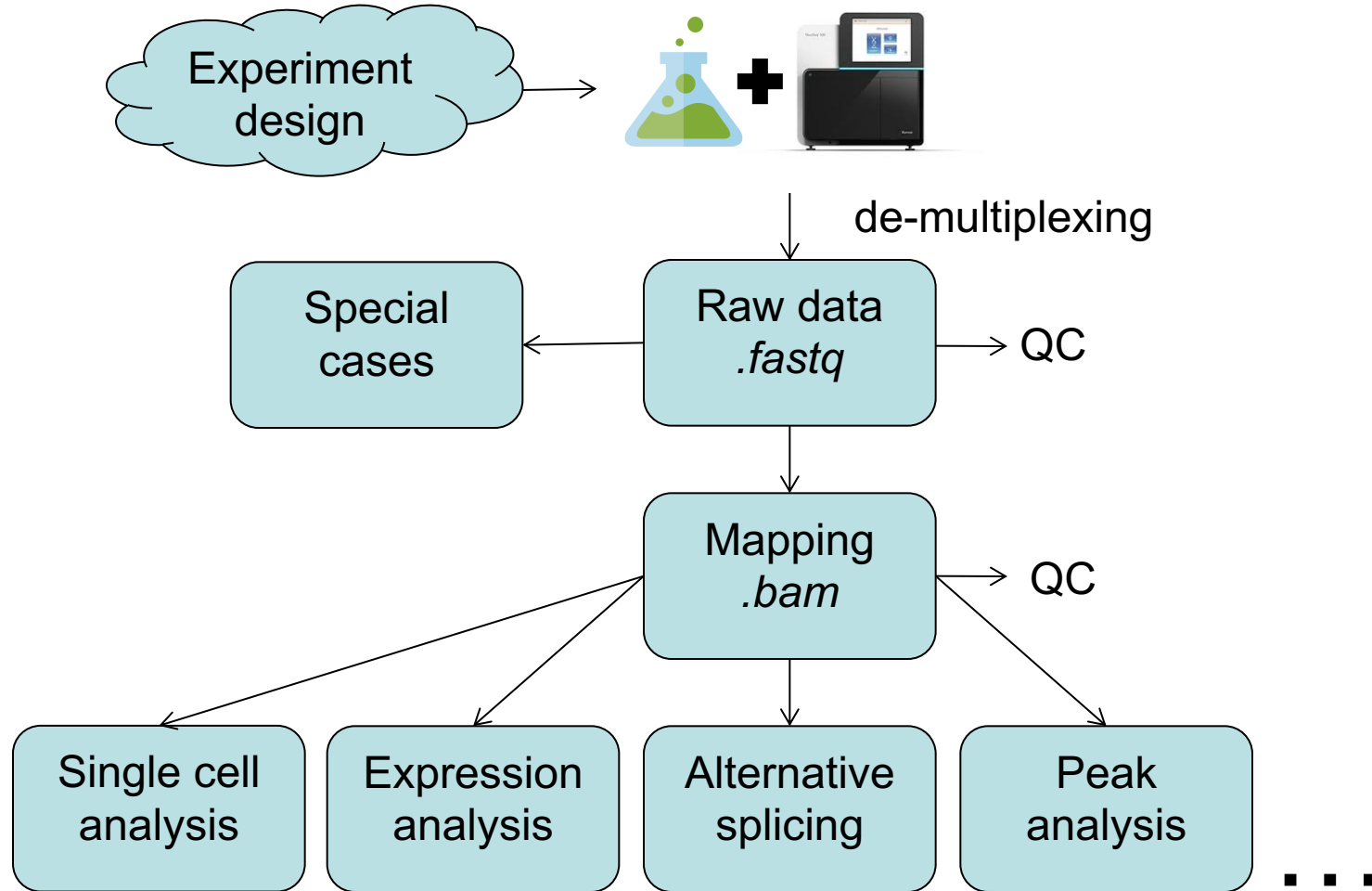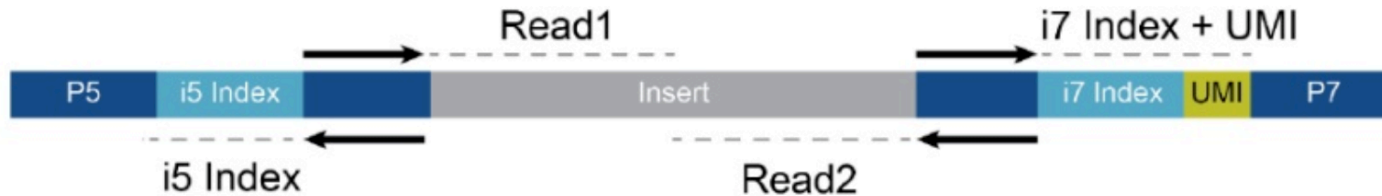CEITEC

# NGS data analysis workflow

# UMI – unique molecular identifiers



- Each molecular fragment gets unique n-base sequence (n ~ 8-12)

- Usage:
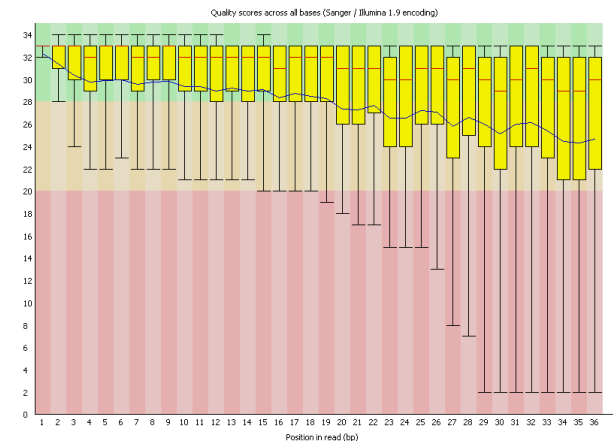  - Mark duplicates

CEITEC

# Raw data - QC

- **Fastq - q stands for quality – coded phred score**

  ```
  CFFFFEFFGCEEGECFGGGGAFF87@E:++6C<++3:,8,33,,:,,,:,,:,,,
  ```

$$Q = -10 \cdot \log_{10} P$$

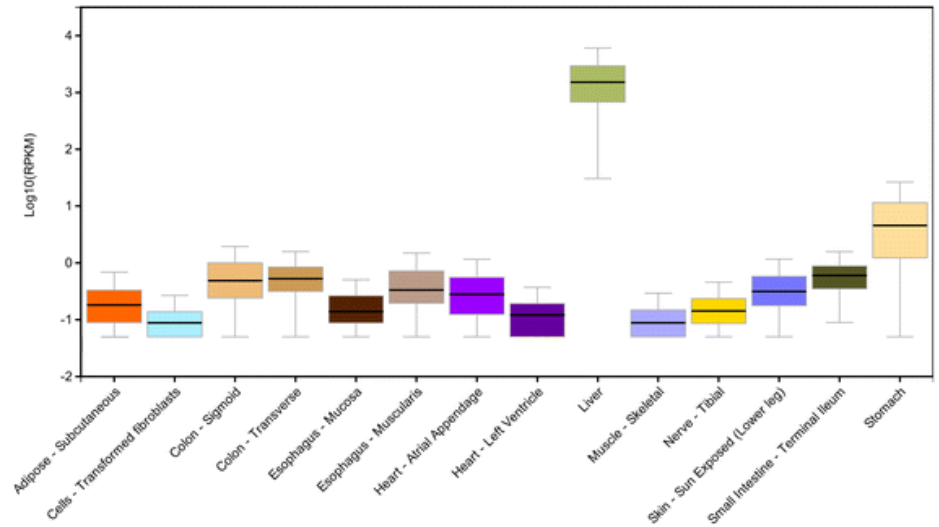| Quality | Error probability |
|---------|-------------------|
| 5 | 31% |
| 10 | 10% |
| 20 | 1% |
| 30 | 0.1% |

- **Very good for early problem detection**
- **Reasonable for trimming and read filtering**
  - **RNA seq  - above phred score 5**

Quality scores across all bases (Sanger / Illumina 1.9 encoding)

CEITEC

# Alignment - QC

- **Per gene coverage**

- **Variability of per gene mapping**

- **Gene counts distribution**

- **rRNA content estimate**

- **Tissue expression check - gtex**

# Alignment - QC

- **QC example – multiQC html**

# NGS data analysis workflow

# Expression analysis - planning

- **3 way balance**
  - **Read depth**
  - **Biological replicates**
  - **Fold change (number of genes) sensitivity**

very sensitive

many BR — high RD

low RD — not sensitive — 1 BR

Table 1

Statistical power to detect differential expression varies with effect size, sequencing depth and number of replicates

| | Replicates per group | | |
|---|---|---|---|
| | 3 | 5 | 10 |
| **Effect size (fold change)** | | | |
| 1.25 | 17 % | 25 % | 44 % |
| 1.5 | 43 % | 64 % | 91 % |
| 2 | 87 % | 98 % | 100 % |
| **Sequencing depth (millions of reads)** | | | |
| 3 | 19 % | 29 % | 52 % |
| 10 | 33 % | 51 % | 80 % |
| 15 | 38 % | 57 % | 85 % |

CEITEC

# Expression analysis - planning

- **Depth**

- **Human ~ 22 000 genes = minimum 20 mil mapped reads**

- **Good 25 mil mapped reads**


- **Mapped reads!**

  - **rRNA removal**

  - **Size selection for sRNA**


- **Technical vs. biological**

  - **Technical only for technique testing**

- **Batch effect**

  - **Sample randomized sequencing**

- **Highly suggested minimum = 4 rep**

CEITEC

# Expression analysis

- **Raw counts**

# Expression analysis

- **Result**

# Expression analysis

- **Result**
  - **normCounts**
    - **rpkm - Reads Per Kilobase of transcript per Million mapped reads**
    - **fpkm - Fragments Per Kilobase of transcript per Million mapped reads**
    - **tpm - Transcripts Per Million (TPM)**
      - **for every 1,000,000 RNA molecules in the RNA-seq sample, x came from this gene/transcript**
  - **log2FoldChange**
  - **pvalue**
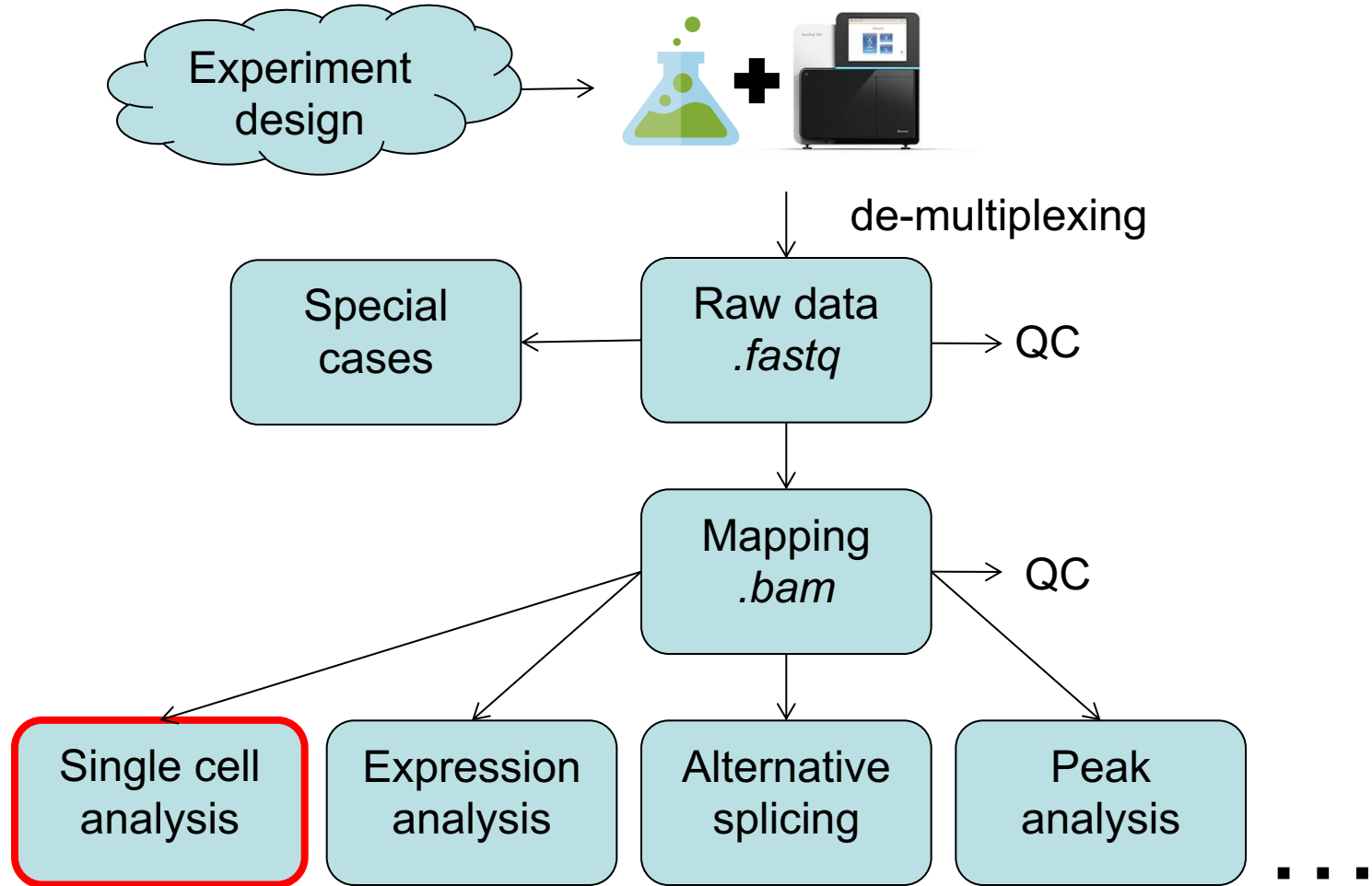  - **padj – pvalue adjusted for multiple testing**

CEITEC

# Expression analysis

- **Report example**

CEITEC

# NGS data analysis workflow

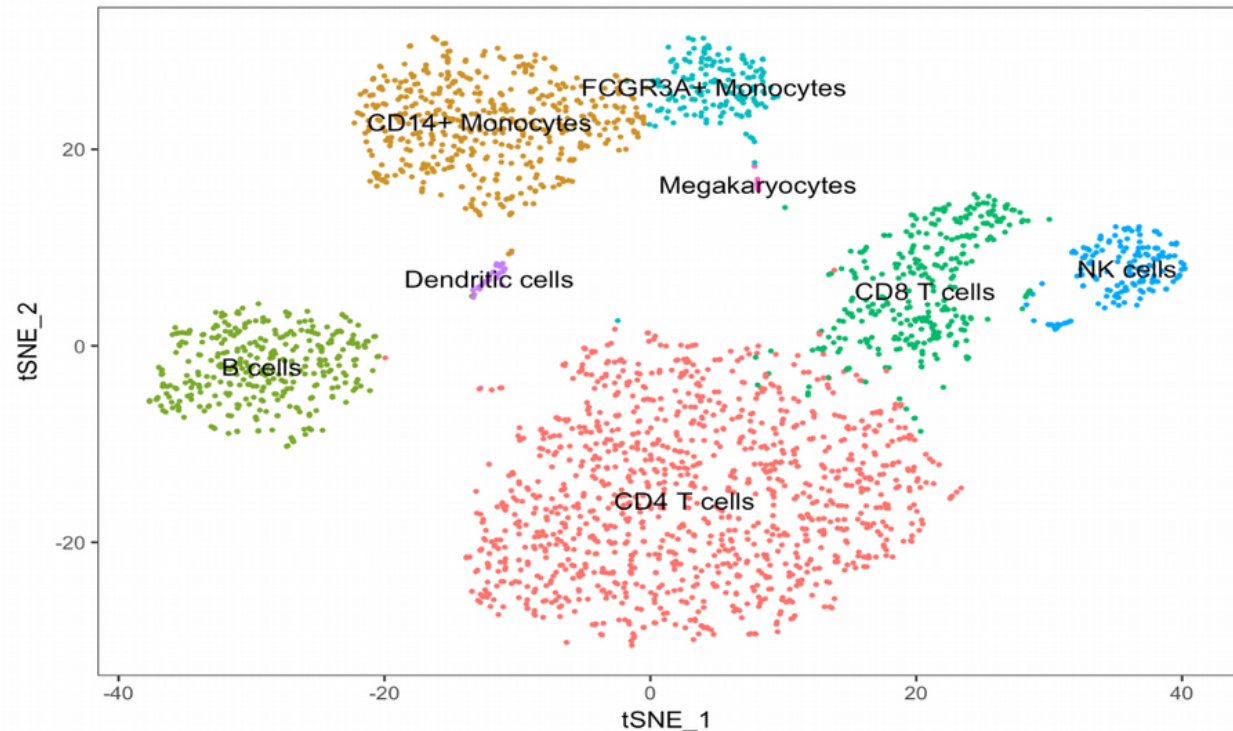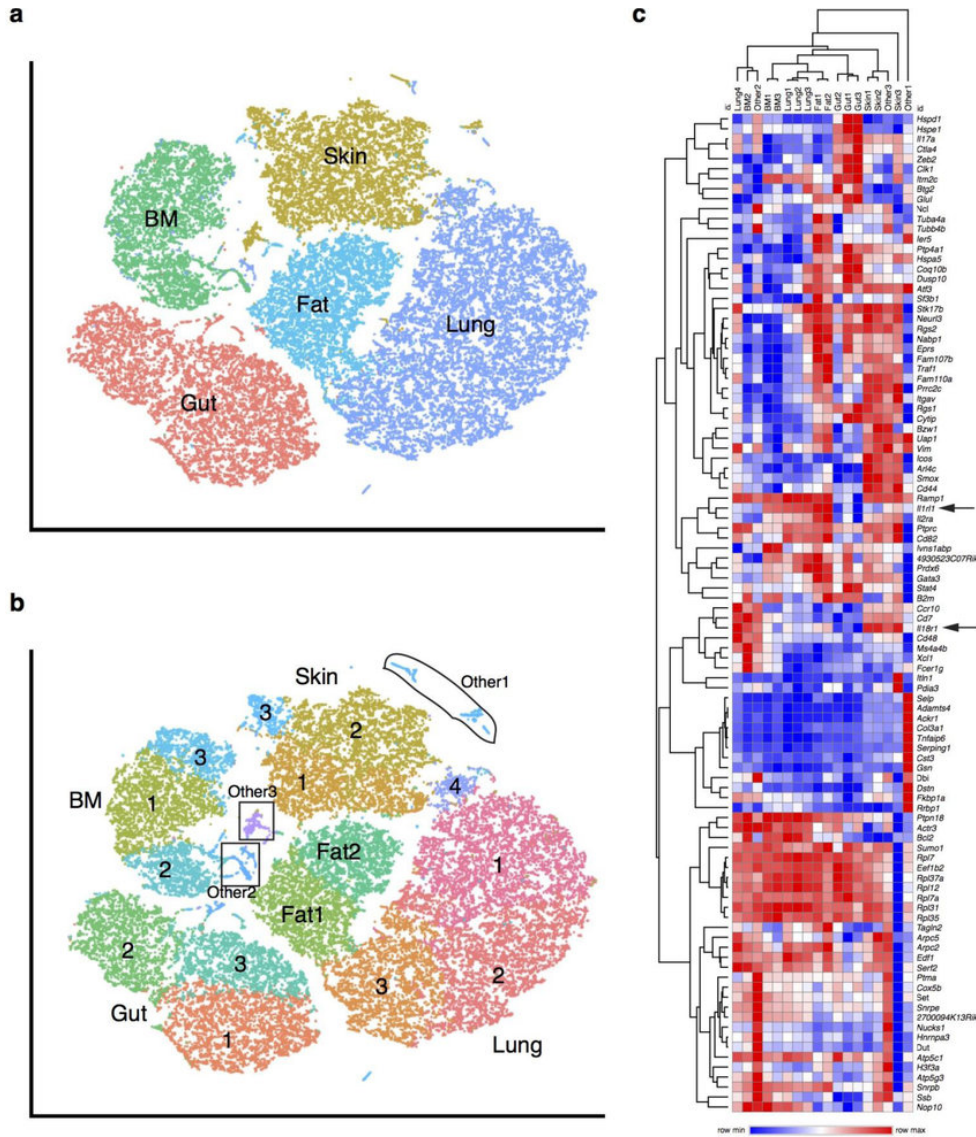# Single cell analysis

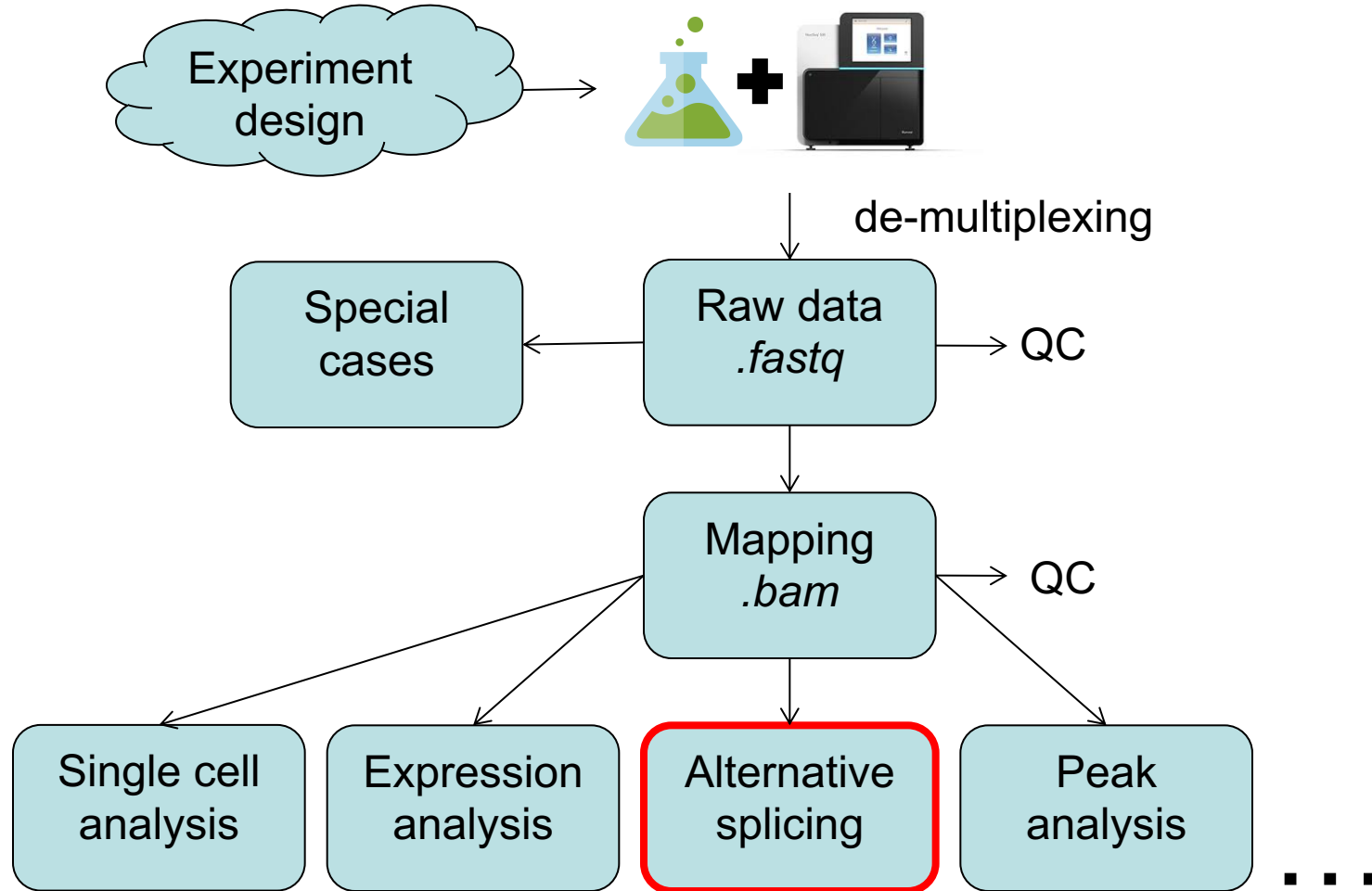- **Cluster cells based on expression**
  - **Cleaning/Filtering step**
  - **Clustering**
  - **Dimension reduction**
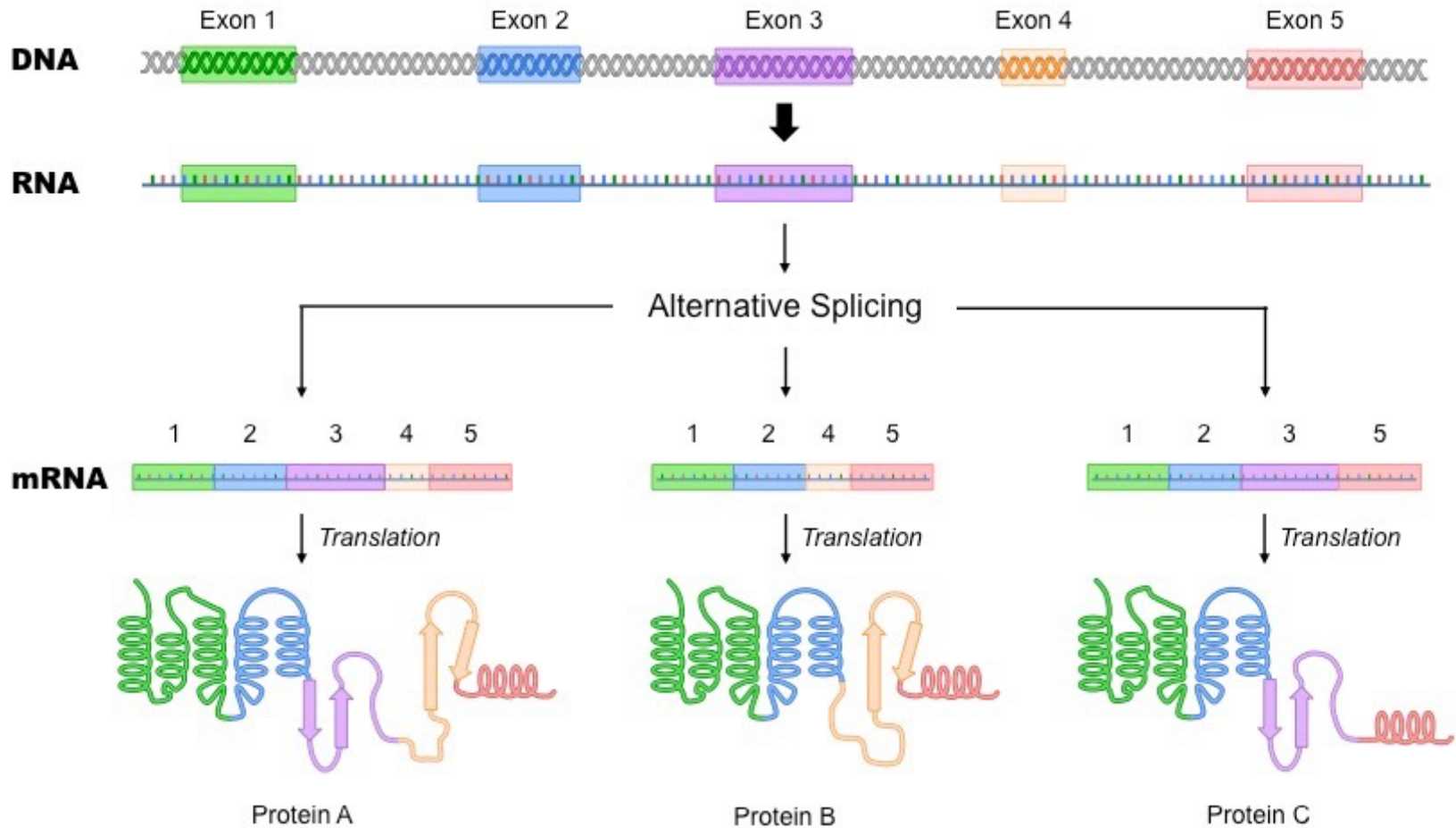    - **PCA**
    - **tSNE**
  - **Visualization**

# Single cell analysis

# NGS data analysis workflow

# Alternative splicing

# Alternative splicing



Exon skipping/inclusion

Alternative 3′ splice sites

Alternative 5′ splice sites

Mutually exclusive exons

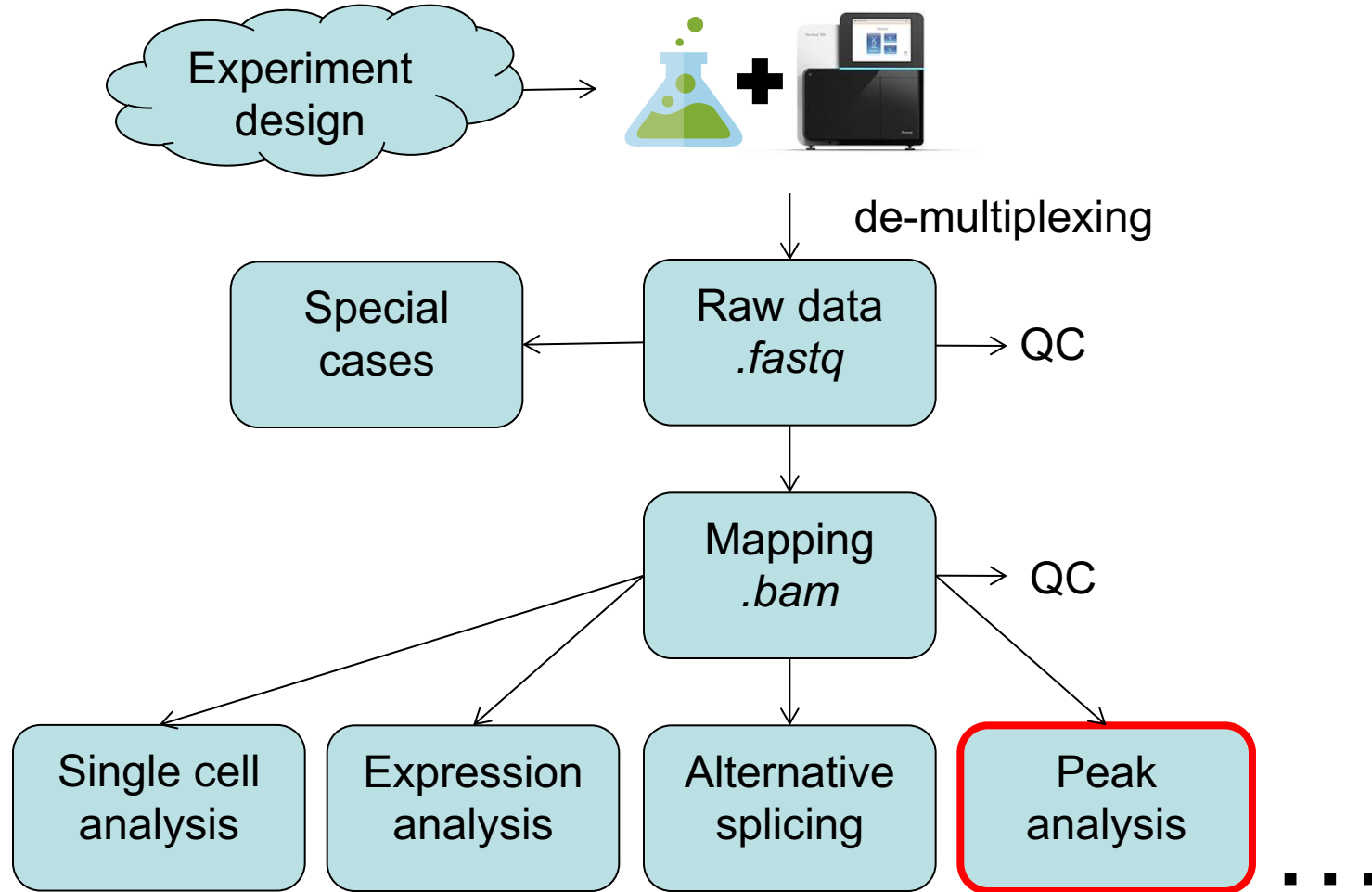Intron retention

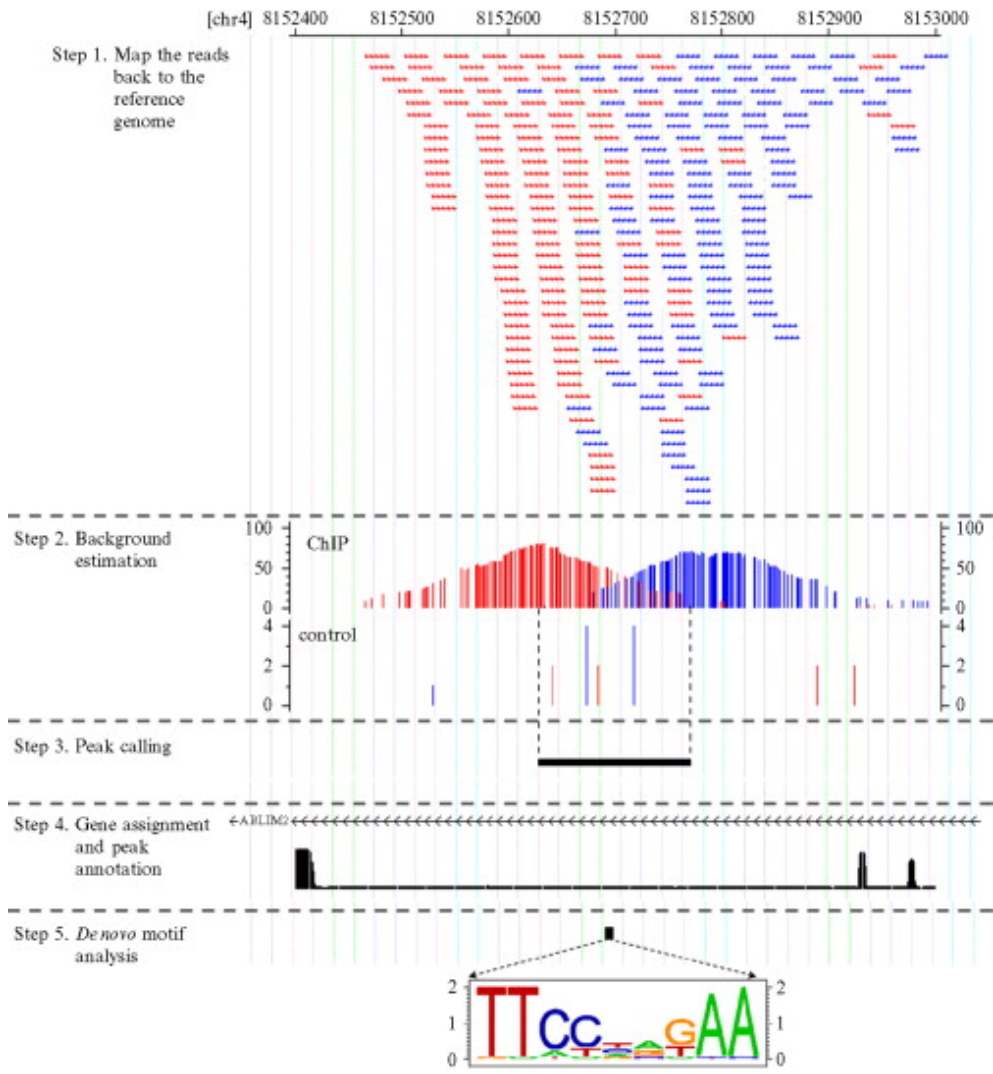Constitutive exon    Alternatively spliced exon

CEITEC

# Alternative splicing
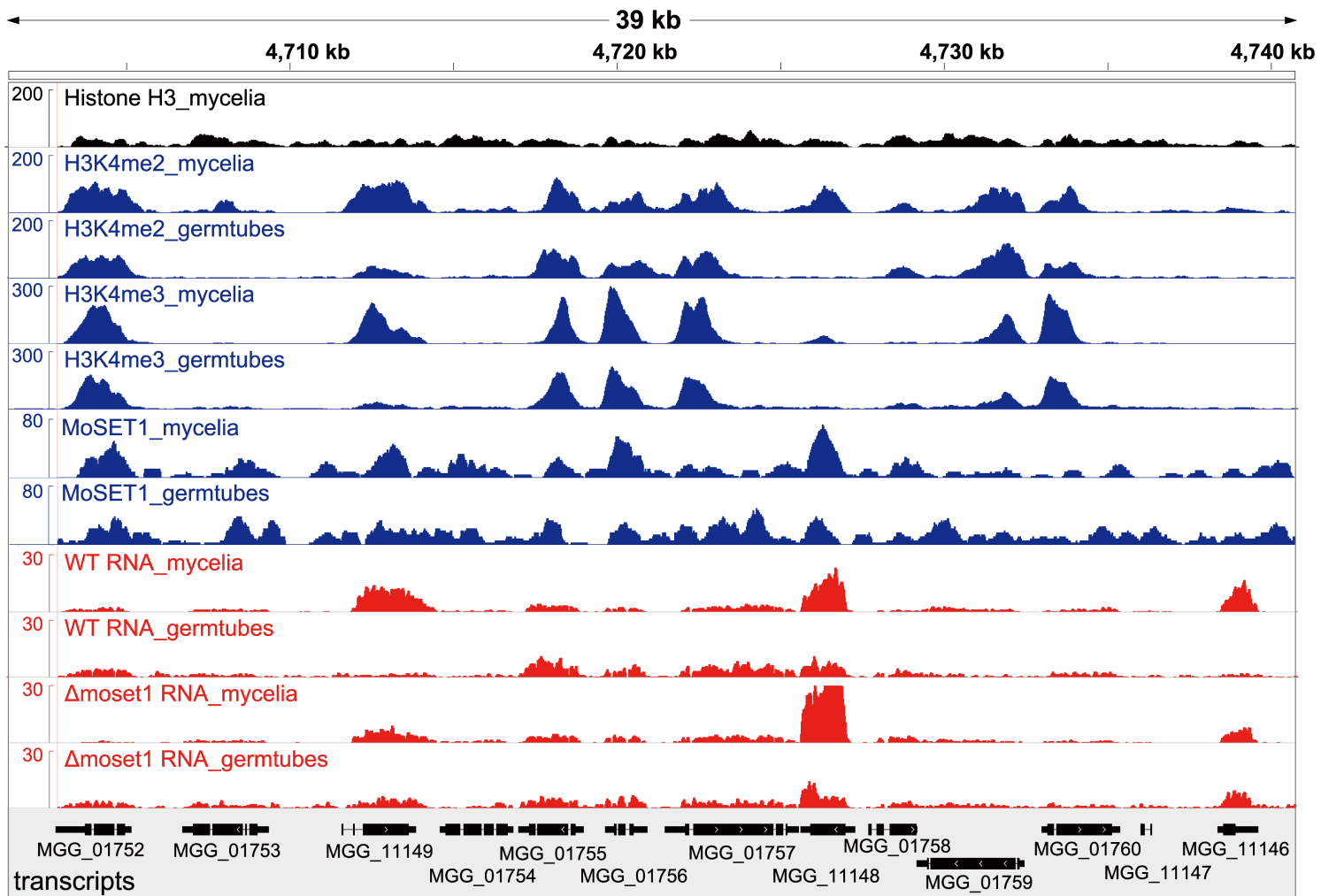
# NGS data analysis workflow

# Peak analysis

# Peak analysis

# Thank you for your attention

CEITEC

Central European Institute of Technology
Masaryk University
Kamenice 753/5
625 00 Brno, Czech Republic

www.ceitec.muni.cz | info@ceitec.muni.cz